

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2002-244894

(43)Date of publication of application : 30.08.2002

(51)Int.Cl.

G06F 12/00
G06F 5/00
G06F 13/00
H03M 7/30
// G06F 17/21

(21)Application number : 2001-027462

(71)Applicant : INTERNATL BUSINESS MACH CORP
<IBM>

(22)Date of filing : 02.02.2001

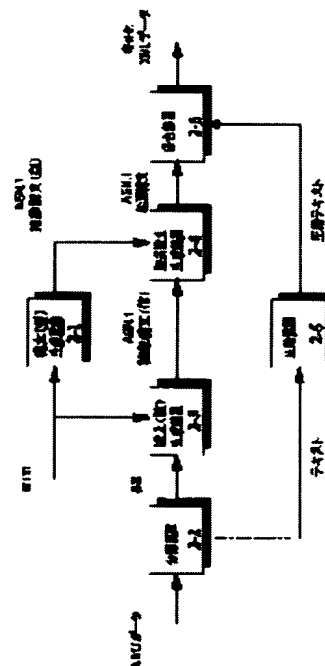
(72)Inventor : IMAMURA TAKESHI

(54) ENCODING METHOD FOR XML DATA, DECODING METHOD FOR ENCODED XML DATA, ENCODING SYSTEM FOR XML DATA, DECODING SYSTEM FOR ENCODED XML DATA, PROGRAM AND RECORDING MEDIUM

(57)Abstract:

PROBLEM TO BE SOLVED: To enhance the efficiency of encoding (compression) of XML data.

SOLUTION: A syntax (type) generator 2-1 is used to convert a DTD into ASN.1 abstract syntax (type), and a separator 2-2 is used to separate text from XML data according to the DTD. A syntax (value) generator 2-3 is next used to convert the XML data (syntax of elements) from which the text is separated into ASN.1 abstract syntax (value) according to the ASN.1 abstract syntax (type). A transfer syntax generator 2-4 is then used to convert the ASN.1 abstract syntax (value) into ASN.1 transfer syntax. A compressor 2-5 is used to compress the separated text. A merger 2-6 is further used to merge the ASN.1 transfer syntax and the compressed text to generate encoded XML data.



LEGAL STATUS

[Date of request for examination] 16.10.2001

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

(19)日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11)特許出願公開番号

特開2002-244894

(P2002-244894A)

(43)公開日 平成14年8月30日(2002.8.30)

(51)Int.Cl. ⁷	識別記号	F I	テーマコード(参考)
G 0 6 F 12/00	5 1 1	G 0 6 F 12/00	5 1 1 A 5 B 0 0 9
	5 4 6		5 4 6 A 5 B 0 8 2
5/00		5/00	H 5 B 0 8 9
13/00	3 5 3	13/00	3 5 3 C 5 J 0 6 4
H 0 3 M 7/30		H 0 3 M 7/30	Z

審査請求 有 請求項の数23 OL (全 17 頁) 最終頁に続く

(21)出願番号 特願2001-27462(P2001-27462)

(22)出願日 平成13年2月2日(2001.2.2)

(71)出願人 390009531

インターナショナル・ビジネス・マシー

ズ・コーポレーション

INTERNATIONAL BUSIN

ESS MASCHINES CORPO

RATION

アメリカ合衆国10504、ニューヨーク州

アーモンク ニュー オーチャード ロー

ド

(74)復代理人 100112520

弁理士 林 茂則 (外3名)

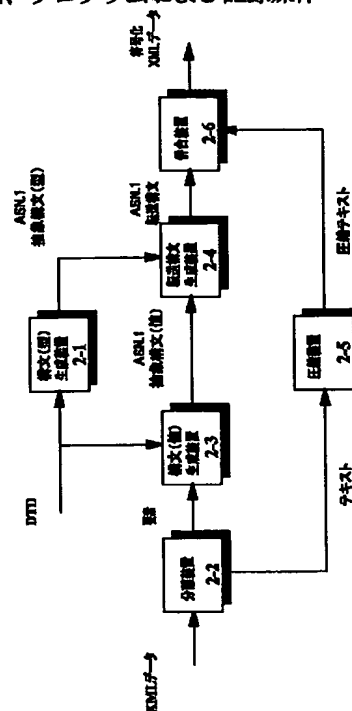
最終頁に続く

(54)【発明の名称】 XMLデータの符号化方法、符号化されたXMLデータの復号化方法、XMLデータの符号化システム、符号化されたXMLデータの復号化システム、プログラムおよび記録媒体

(57)【要約】

【課題】 XMLデータの符号化(圧縮)効率を高める。

【解決手段】 構文(型)生成装置2-1を用いてDTDをASN. 1抽象構文(型)に変換し、分離装置2-2を用いて前記DTDに従うXMLデータからテキストを分離する。その後、構文(値)生成装置2-3を用いてテキストが分離された後のXMLデータ(要素の構文)を、前記ASN. 1抽象構文(型)に従うASN. 1抽象構文(値)に変換する。その後、転送構文生成装置2-4を用いてASN. 1抽象構文(値)をASN. 1転送構文に変換する。一方、圧縮装置2-5を用いて分離されたテキストを圧縮する。さらに、併合装置2-6を用いてASN. 1転送構文と圧縮テキストを併合し符号化XMLデータを生成する。



【特許請求の範囲】

【請求項1】 XMLデータの文法を定義する文法定義をASN. 1抽象構文(型)に変換するステップと、前記XMLデータを、要素の内容(テキスト)と構造(構造を含む要素名)とに分離するステップと、前記構造を、前記ASN. 1抽象構文(型)に従うASN. 1抽象構文(値)に変換するステップと、前記ASN. 1抽象構文(値)をASN. 1転送構文に変換するステップと、前記要素の内容を圧縮するステップと、圧縮された前記要素の内容と前記ASN. 1転送構文とを併合するステップと、を含むXMLデータの符号化方法。

【請求項2】 前記文法定義に属性、処理命令その他前記要素以外の文法定義項目を有し、前記要素以外の文法定義項目を特別要素として前記要素に含めるように前記文法定義を他の文法定義に変換するステップと、前記他の文法定義に従うように、前記XMLデータを他のXMLデータに変換するステップと、をさらに有する請求項1記載のXMLデータの符号化方法。

【請求項3】 前記文法定義はDTDであり、前記DTDの要素内容には、「|」、「?」、「*」または「+」から選択される単一または複数のオペレータが含まれ、または前記オペレータを含まず、前記ASN. 1抽象構文(型)において、前記「|」オペレータは、sequence型で表現し、前記「|」オペレータは、choice型で表現し、前記「?」オペレータは、sequence型とキーワード「OPTIONAL」の組み合わせで表現し、前記「*」オペレータは、sequence-of型で表現し、前記「+」オペレータは、サイズが制限されたsequence-of型で表現し、前記オペレータが無いときには、defined型で表現する、請求項1または2記載のXMLデータの符号化方法。

【請求項4】 前記文法定義を他の文法定義に変換するステップにおいて、前記文法定義に含まれる属性を、一意に決定できる属性要素として表現し、前記属性の親要素の子要素として扱い、前記属性の属性値は、「CDATA」として前記属性要素の子要素として扱い、前記属性が「REQUIRED」属性の場合、前記属性要素を要素で表現し、前記属性が「IMPLIED」属性の場合、または、前記属性にデフォルト値が定義されている場合、前記属性

要素を「?」オペレータが適用された要素で表現し、前記XMLデータを他のXMLデータに変換するステップにおいて、

前記XMLデータの要素に含まれる属性を、一意に決定できる属性要素として表現し、前記属性の親要素の子要素として扱う(ただし前記属性にデフォルト値が定義され、かつ前記属性の属性値が前記デフォルト値と一致する場合を除く)、

請求項3記載のXMLデータの符号化方法。

10 【請求項5】 前記ASN. 1抽象構文(値)を前記ASN. 1転送構文に変換するステップにおいて、PER規則を用いる請求項1または2記載のXMLデータの符号化方法。

【請求項6】 XMLデータの文法を定義する文法定義をASN. 1抽象構文(型)に変換するステップと、符号化されたXMLデータを、ASN. 1転送構文と圧縮された要素の内容(テキスト)とに分離するステップと、

前記ASN. 1転送構文を、前記ASN. 1抽象構文(型)に従うASN. 1抽象構文(値)に変換するステップと、

前記ASN. 1抽象構文(値)を、前記文法定義に従うXMLデータの構造(構造を含む要素名)に変換するステップと、

前記圧縮された要素の内容を解凍するステップと、

前記解凍された要素の内容と前記XMLデータの構造とを併合するステップと、

を含む符号化されたXMLデータの復号化方法。

【請求項7】 前記文法定義に属性、処理命令その他前記要素以外の文法定義項目を有し、前記要素以外の文法定義項目を特別要素として前記要素に含めるように前記文法定義を他の文法定義に変換するステップと、

前記文法定義に従うように、前記復号化されたXMLデータを他のXMLデータに変換するステップと、

をさらに有する請求項6記載の符号化されたXMLデータの復号化方法。

【請求項8】 前記文法定義はDTDであり、

前記DTDの要素内容には、「|」、「?」、「*」または「+」から選択される単一または複数のオペレータが含まれ、または前記オペレータを含まず、前記ASN. 1抽象構文(型)において、前記「|」オペレータは、sequence型で表現し、

前記「|」オペレータは、choice型で表現し、

前記「?」オペレータは、sequence型とキーワード「OPTIONAL」の組み合わせで表現し、

前記「*」オペレータは、sequence-of型で表現し、

50 前記「+」オペレータは、サイズが制限されたsequ

3

ence-of型で表現し、

前記オペレータが無いときには、defined型で表現する、

請求項6または7記載の符号化されたXMLデータの復号化方法。

【請求項9】 前記文法定義を他の文法定義に変換するステップにおいて、

前記文法定義に含まれる属性を、一意に決定できる属性要素として表現し前記属性の親要素の子要素として扱い、

前記属性の属性値は、「CDATA」として前記属性要素の子要素として扱い、

前記属性が「REQUIRED」属性の場合、前記属性要素を要素で表現し、

前記属性が「IMPLIED」属性の場合、または、前記属性にデフォルト値が定義されている場合、前記属性要素を「？」オペレータが適用された要素で表現し、

前記復号化されたXMLデータを他のXMLデータに変換するステップにおいて、

前記復号化されたXMLデータの要素に子要素として含まれる前記属性要素を、前記要素の属性またはその属性値に変換する、

請求項8記載の符号化されたXMLデータの復号化方法。

【請求項10】 前記ASN.1転送構文を前記ASN.1抽象構文(値)に変換するステップにおいて、PER規則を用いる請求項6または7記載の符号化されたXMLデータの復号化方法。

【請求項11】 XMLデータの文法を定義する文法定義をASN.1抽象構文(型)に変換する手段と、

前記XMLデータを、要素の内容(テキスト)と構造(構造を含む要素名)とに分離する手段と、

前記構造を、前記ASN.1抽象構文(型)に従うASN.1抽象構文(値)に変換する手段と、

前記ASN.1抽象構文(値)をASN.1転送構文に変換する手段と、

前記要素の内容を圧縮する手段と、

圧縮された前記要素の内容と前記ASN.1転送構文とを併合する手段と、

を含むXMLデータの符号化システム。

【請求項12】 前記文法定義に属性、処理命令その他前記要素以外の文法定義項目を有し、

前記要素以外の文法定義項目を特別要素として前記要素に含めるように前記文法定義を他の文法定義に変換する手段と、

前記他の文法定義に従うように、前記XMLデータを他のXMLデータに変換する手段と、

をさらに有する請求項11記載のXMLデータの符号化システム。

【請求項13】 前記文法定義はDTDであり、

4

前記DTDの要素内容には、「|」、「?」、「*」または「+」から選択される単一または複数のオ

ペレータが含まれ、または前記オペレータを含まず、

前記ASN.1抽象構文(型)において、

前記「|」オペレータは、sequence型で表現し、

前記「|」オペレータは、choice型で表現し、

前記「?」オペレータは、sequence型とキーワード「OPTIONAL」の組み合わせで表現し、

10 前記「*」オペレータは、sequence-of型で表現し、

前記「+」オペレータは、サイズが制限されたsequence-of型で表現し、

前記オペレータが無いときには、defined型で表現する、

請求項11または12記載のXMLデータの符号化システム。

【請求項14】 前記文法定義を他の文法定義に変換する手段において、

20 前記文法定義に含まれる属性を、一意に決定できる属性要素として表現し、前記属性の親要素の子要素として扱い、

前記属性の属性値は、「CDATA」として前記属性要素の子要素として扱い、

前記属性が「REQUIRED」属性の場合、前記属性要素を要素で表現し、

前記属性が「IMPLIED」属性の場合、または、前記属性にデフォルト値が定義されている場合、前記属性要素を「？」オペレータが適用された要素で表現し、

30 前記XMLデータを他のXMLデータに変換する手段において、

前記XMLデータの要素に含まれる属性を、一意に決定できる属性要素として表現し、前記属性の親要素の子要素として扱う(ただし前記属性にデフォルト値が定義され、かつ前記属性の属性値が前記デフォルト値と一致する場合を除く)、

請求項13記載のXMLデータの符号化システム。

【請求項15】 前記ASN.1抽象構文(値)を前記ASN.1転送構文に変換する手段において、PER規則を用いる請求項11または12記載のXMLデータの

40 符号化システム。

【請求項16】 XMLデータの文法を定義する文法定義をASN.1抽象構文(型)に変換する手段と、

符号化されたXMLデータを、ASN.1転送構文と圧縮された要素の内容(テキスト)とに分離する手段と、

前記ASN.1転送構文を、前記ASN.1抽象構文(型)に従うASN.1抽象構文(値)に変換する手段と、

前記ASN.1抽象構文(値)を、前記文法定義に従うXMLデータの構造(構造を含む要素名)に変換する手

50

段と、

前記圧縮された要素の内容を解凍する手段と、
前記解凍された要素の内容と前記XMLデータの構造とを併合する手段と、
を含む符号化されたXMLデータの復号化システム。

【請求項17】 前記文法定義に属性、処理命令その他前記要素以外の文法定義項目を有し、
前記要素以外の文法定義項目を特別要素として前記要素に含めるように前記文法定義を他の文法定義に変換する手段と、

前記文法定義に従うように、前記復号化されたXMLデータを他のXMLデータに変換する手段と、
をさらに有する請求項16記載の符号化されたXMLデータの復号化システム。

【請求項18】 前記文法定義はDTDであり、
前記DTDの要素内容には、「,」、「|」、「?」、「*」または「+」から選択される単一または複数のオペレータが含まれ、または前記オペレータを含まず、
前記ASN.1抽象構文(型)において、
前記「,」オペレータは、sequence型で表現し、
前記「|」オペレータは、choice型で表現し、
前記「?」オペレータは、sequence型とキーワード「OPTIONAL」の組み合わせで表現し、
前記「*」オペレータは、sequence-of型で表現し、
前記「+」オペレータは、サイズが制限されたsequence-of型で表現し、
前記オペレータが無いときには、defined型で表現する、
請求項16または17記載の符号化されたXMLデータの復号化システム。

【請求項19】 前記文法定義を他の文法定義に変換する手段において、
前記文法定義に含まれる属性を、一意に決定できる属性要素として表現し、前記属性の親要素の子要素として扱い、
前記属性の属性値は、「CDATA」として前記属性要素の子要素として扱い、
前記属性が「REQUIRED」属性の場合、前記属性要素を要素で表現し、
前記属性が「IMPLIED」属性の場合、または、前記属性にデフォルト値が定義されている場合、前記属性要素を「?」オペレータが適用された要素で表現し、
前記復号化されたXMLデータを他のXMLデータに変換する手段において、
前記復号化されたXMLデータの要素に子要素として含まれる前記属性要素を、前記要素の属性またはその属性値に変換する、
請求項18記載の符号化されたXMLデータの復号化シ

ステム。

【請求項20】 前記ASN.1転送構文を前記ASN.1抽象構文(値)に変換する手段において、PER規則を用いる請求項16または17記載の符号化されたXMLデータの復号化システム。

【請求項21】 コンピュータが実行可能なプログラムであって、コンピュータに、
XMLデータの文法を定義する文法定義をASN.1抽象構文(型)に変換する機能と、

10 前記XMLデータを、要素の内容(テキスト)と構造(構造を含む要素名)とに分離する機能と、
前記構造を、前記ASN.1抽象構文(型)に従うASN.1抽象構文(値)に変換する機能と、

前記ASN.1抽象構文(値)をASN.1転送構文に変換する機能と、

前記要素の内容を圧縮する機能と、

圧縮された前記要素の内容と前記ASN.1転送構文とを併合する機能と、
を実現させるためのプログラム。

20 【請求項22】 コンピュータが実行可能なプログラムであって、コンピュータに、
XMLデータの文法を定義する文法定義をASN.1抽象構文(型)に変換する機能と、
符号化されたXMLデータを、ASN.1転送構文と圧縮された要素の内容(テキスト)とに分離する機能と、
前記ASN.1転送構文を、前記ASN.1抽象構文(型)に従うASN.1抽象構文(値)に変換する機能と、

前記ASN.1抽象構文(値)を、前記文法定義に従うXMLデータの構造(構造を含む要素名)に変換する機能と、
30 前記圧縮された要素の内容を解凍する機能と、
前記解凍された要素の内容と前記XMLデータの構造とを併合する機能と、
を実現させるためのプログラム。

【請求項23】 コンピュータ読取可能な記録媒体であって、
コンピュータに、XMLデータの文法を定義する文法定義をASN.1抽象構文(型)に変換する機能と、前記XMLデータを、要素の内容(テキスト)と構造(構造を含む要素名)とに分離する機能と、前記構造を、前記ASN.1抽象構文(型)に従うASN.1抽象構文(値)に変換する機能と、前記ASN.1抽象構文(値)をASN.1転送構文に変換する機能と、前記要素の内容を圧縮する機能と、圧縮された前記要素の内容と前記ASN.1転送構文とを併合する機能と、を実現させるための第1プログラム、または、
コンピュータに、XMLデータの文法を定義する文法定義をASN.1抽象構文(型)に変換する機能と、符号化されたXMLデータを、ASN.1転送構文と圧縮さ

れた要素の内容(テキスト)とに分離する機能と、前記ASN.1転送構文を、前記ASN.1抽象構文(型)に従うASN.1抽象構文(値)に変換する機能と、前記ASN.1抽象構文(値)を、前記文法定義に従うXMLデータの構造(構造を含む要素名)に変換する機能と、前記圧縮された要素の内容を解凍する機能と、前記解凍された要素の内容と前記XMLデータの構造とを併合する機能と、を実現させるための第2プログラム、の何れかのプログラムが記録された記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、データの符号化および復号化技術に関し、特にXML(extensible markup language)データの符号化における圧縮率の向上に適用して有効な技術に関する。

【0002】

【従来の技術】近年、インターネット上でのデータ表現の手段としてXMLが目まぐるしく注目を浴びている。XMLは拡張可能なメタ言語であり、ユーザが独自に文法を規定することができる。また、各要素に論理的な意味を持たせることが可能であり、HTML(hypertext markup language)に比較してデータ処理が大幅に容易になる。このため、電子商取引等インターネットで交換される構造化文書の標準的な表現方式として期待されている。なお、XMLについては、「W3C. Extensible Markup Language (XML) 1.0, 1998. <http://www.w3.org/TR/REC-xml>」に詳細が記述されている。

【0003】XMLデータは文字で記述されるため可読性が高い利点がある。その反面、冗長性も高いという難点がある。すなわち要素の意味は主に開始タグ内に記述され、その意味内容は文字で記述されている内容を参照すれば容易に把握できる。しかし、全てが文字で記述されるため、全体の文字数が増加し、XMLデータ全体の情報量(文字数)が多くなる。文字数が多ければ、たとえばストレージに格納したり、あるいはネットワークで転送する際に、記録容量あるいは転送量が多くなり物理的、時間的コストの上昇を来す。そこで、XMLデータを短い符号に符号化(または圧縮)できれば都合がよい。

【0004】データの圧縮法には各種の手法が知られている。たとえば、ランレングス符号、Huffman符号、算術符号、LZ77等である。これら符号化の手法については、たとえば、「Huffman, D.A. "A method for the construction of minimum-redundancy codes" Proc. of the IRE September, 1952」、「Mark Nelson and Jean Loup Gailly "The Data Compression Book", Second Edition. M&TBooks 1996」、「Jacob Ziv and Abraham Lempel. "A universal algorithm for sequential data compression" IEEE Transactions on Information Theory May, 1977」に詳しい。

【0005】しかしながら、これら圧縮手法はXMLに特化されたものではなく、XMLデータに適用した場合に必ずしも圧縮効率が高いものではない。XMLデータに特化した圧縮手法には、たとえば「D. Suciu and H. Liefke. XMill: an Efficient Compressor for XML Data, 1999. <http://www.research.att.com/sw/tools/xmill/>」に記載のXMill、「XML Solutions Corp. XML Zip, 1999. <http://www.xmls.com/products/xmlzip/xmlzip.html>」に記載のXMLZip、あるいは「井川甲作 著、東京工業大学工学部情報工学科卒業論文「DTDを用いたXML文書圧縮アルゴリズムに関する研究」、平成12年2月」に記載のXCompなどがある。

【0006】XMillは、XMLデータから要素ごとのコンテンツ(テキスト)部分を抽出する。この抽出された部分をコンテナと呼ぶ。そして構造部分を数字で符号化し、テキスト部分はコンテナごとにLZ77などの方法で圧縮する。基本的にはパラメータ等の情報を必要とせず、アプリケーションのみで圧縮が可能である。必要であればパラメータ等を指定してコンテナごとの圧縮方法を指定し、圧縮効率を高めることが可能である。また、Cで実装されるため圧縮速度が速いという特徴をもつ。

【0007】XMLZipは、ルート要素からの深さを指定し、指定部分をドキュメント要素から分割し、残りをZIPで圧縮する。ルート要素部分は符号化されず直接操作することが可能になる。使用しない部分を圧縮することで文書へのアクセスを迅速に行える。ただし、圧縮効率はXMillに比較して低い。

【0008】XCompは、XMLデータの構造部分のうち、DTD(Document Type Definition)から一意に決定される部分は符号化せず、一意に決定できない部分についての構造部分のみを圧縮する。テキスト部分はXMillと同じ方法で圧縮する。すなわち、以下の手順で圧縮を行う。(1)構造とコンテンツを分離する。

(2)DTDからプッシュダウンオートマトン(PDA)を生成する。(3)生成されたPDAを用いて構造部分を符号化する符号化トランスデューサを生成する。

(4)符号化トランスデューサの各ノードに割り振られた数字をオートマトンを連鎖的に遷移することにより出力し、構造を符号化する。(5)得られた構造符号と要素ごとのコンテンツをLZ77等で圧縮し、圧縮されたXML文書を出力する。

【0009】上記XML文書に特化された圧縮手法のうち、XCompは構造部分の一部を符号化しないので、かなりよい圧縮効率を達成する。

【0010】

【発明が解決しようとする課題】上記の通り、XCompはXML文書の圧縮手法として相対的に優れる。しか

し、本発明者の検討によればXCompにおいても、XMLデータがある特定の構造をしている場合に圧縮効率が悪くなるという課題がある。すなわち、要素に「？」オペレータや「*」オペレータ（「+」オペレータを含む）が適用された場合である。

【0011】「？」オペレータは、ある要素の子要素が0回もしくは1回出現する場合に、要素の宣言文において子要素に付されるオペレータである。XCompにおいて「？」オペレータが出現した場合、「？」オペレータは、ある状態から別の状態に移移する複数の選択肢で表現される。選択肢にはインデックスが付与される。XCompの実行時には、選択された選択肢に付与されたインデックスが出力される。選択肢の数は、「？」オペレータの後に「？」オペレータなどがいくつ連続するかによる。例えば、「？」オペレータがn個連続する場合に、最初の「？」オペレータに対する選択肢はn+1個となる。従って、インデックスはn+1種類必要であり、その1つを表現するには $O(\log n)$ ビット必要である。「？」オペレータが適用された要素がすべて存在すると、インデックスが複数羅列される。その1つを表現するのに $O(\log n)$ ビット必要であるから、それ全体を表現するには $O(n \log n)$ ビット必要である。

【0012】「*」オペレータは、ある要素の子要素が0回以上出現する場合に、要素の宣言文において子要素に付されるオペレータであり、「+」オペレータは、ある要素の子要素が1回以上出現する場合に、要素の宣言文において子要素に付されるオペレータである。XCompにおいて「*」オペレータ（または「+」オペレータ）が出現した場合、「*（または+）」オペレータは、同じ状態から別の状態に移移する2つの選択肢で表現される。選択肢には、インデックスが付与される。XCompの実行時には、選択された選択肢に付与されたインデックスが出力される。「*」オペレータが適用された要素が複数存在すると、同じインデックスが複数羅列される。インデックスの数は、存在する要素の数に比例する。従って、インデックス全体を表現するのに、要素の数をnとすると、 $O(n)$ ビット必要である。

【0013】すなわち、XCompではDTDから一意的に決まらない構造のうち特定のものについて、それらを表現する符号のビット数が大きくなって、必ずしも十分な圧縮効率を達成できない問題がある。

【0014】本発明の目的は、より圧縮効率の高いXMLデータ（XML文書）の符号化方式を提案し、それを実現する方法およびシステムを提供することにある。

【0015】

【課題を解決するための手段】より高いXMLデータの圧縮効率を達成するために、本発明では、「ITU-T. X.680 - Abstract Syntax Notation One (ASN.1): Specification of basic notation, 1997. <http://www.itu.int>

/itodoc/itu-t/rec/x/x500up/x680.html」に記載のASN.1を利用する。すなわち、符号化対象のXMLデータを構造とコンテンツ（テキスト）に分離する。一方、符号化対象のXMLデータの文法をASN.1抽象構文（型）に変換する。分離したコンテンツは要素ごとあるいはまとめてXMi11と同様に圧縮する。また、分離した構造（要素）をASN.1抽象構文（型）に従うASN.1抽象構文（値）に変換する。次に、ASN.1抽象構文（値）をASN.1で規定されている符号化規則を用いてASN.1転送構文に変換する。ASN.1で規定される符号化規則には、BER、DER、PER等があるが、特に符号化効率の観点からPERを用いることが好ましい。そして符号化された要素（構文）であるASN.1転送構文と圧縮されたコンテンツ（テキスト）を併合し、符号化されたXMLデータを生成する。なお、PERについては、「ITU-T. X.691-ASN.1 encoding rules: Specification of Packed Encoding Rules (PER), 1997. <http://www.itu.int/itodoc/itu-t/rec/x/x500up/x691.html>」に詳しく記載されている。

【0016】このような符号化手法を用いることにより、XCompに比較して圧縮効率の高い符号化が行える。たとえば「？」オペレータがn個連続して出現した場合の構造の符号化では、本発明を用いることにより $O(n)$ ビットでそれを表現できる。また、「*」オペレータが出現した場合の構造の符号化では、本発明を用いることによりオーダとしてはXCompと同じであるが、実質的には少ないビット数でそれを表現できる。そして符号化（圧縮）されたXMLデータを通信あるいはストレージに用いることにより通信負荷を低減し、またストレージの容量を節約できる。

【0017】なお、符号化されたXMLデータの復号化は、前記と逆の処理を行うことにより実現できる。また、本発明の符号化あるいは復号化の方法は、システムとして把握することも可能であり、さらに、前記方法によって実現される機能をコンピュータに実現させるためのプログラムとして把握することも可能である。

【0018】

【発明の実施の形態】以下、本発明の実施の形態を図面を用いて詳細に説明する。ただし、本発明は多くの異なる態様で実施することが可能であり、本実施の形態の記載内容に限定して解釈すべきではない。なお、実施の形態の全体を通して同じ要素には同じ番号を付するものとする。

【0019】以下の実施の形態では、主に方法またはシステムについて説明するが、当業者であれば明らかとなり、本発明はコンピュータで使用可能なプログラムとしても実施できる。したがって、本発明は、ハードウェアとしての実施形態、ソフトウェアとしての実施形態またはソフトウェアとハードウェアとの組合せの実施形態をとることができる。プログラムは、ハードディスク、

CD-ROM、光記憶装置または磁気記憶装置等の任意のコンピュータ可読媒体に記録できる。

【0020】また以下の実施の形態では、一般的なコンピュータシステムを用いることができる。実施の形態で用いることができるコンピュータシステムは、中央演算処理装置（CPU）、主記憶装置（メインメモリ：RAM）、不揮発性記憶装置（ROM）、コプロセッサ、画像アクセラレータ、キャッシュメモリ、入出力制御装置（I/O）等、一般的にコンピュータシステムに備えられるハードウェア資源を備える。また、ハードディスク装置等の外部記憶装置、インターネット等のネットワークに接続可能な通信手段を備えることができる。コンピュータシステムには、パーソナルコンピュータ、ワークステーション、メインフレームコンピュータ等各種のコンピュータが含まれる。

【0021】（実施の形態1）

1. 前提条件

本実施の形態を説明するにあたり、本実施の形態における前提となる条件を述べる。

（1）XMLデータは、文法で構造が定義される。また、文法は、XMLデータを符号化する側と復号化する側で共有される。たとえば文法定義ファイルを外部データとして所定のIPアドレスに記録し、符号化あるいは復号化する際にIPアドレスを用いてこの外部ファイルを参照できる。なお、本実施の形態では文法定義としてDTDを用いるが、これに限られない。例えば、XML SchemaやRELAXなどを用いることができる。

（2）XMLデータは、要素とテキストだけから構成される。その他のもの（例えば、属性や処理命令）は、別の方法で管理される。例えば、属性や処理命令を特別な要素として表現し、XMLデータの中に埋め込むことができる。この際、文法も変更する必要がある。また、XMLデータから分離して、XPathなどと共に別に格納することも考えられる。

（3）テキストは、XMLなどの方で圧縮される。このような前提条件の下に本実施の形態のXMLデータの符号化および復号化が行われることを説明する。

【0022】2. 符号化手順

2.1 システム構成と手順の概要

図1は本実施の形態の符号化システムの一例をその機能について示したブロック図である。本実施の形態の符号化システムは、構文（型）生成装置2-1、分離装置2-2、構文（値）生成装置2-3、転送構文生成装置2-4、圧縮装置2-5、併合装置2-6を有する。

【0023】構文（型）生成装置2-1は、DTDからASN.1抽象構文（型）を生成し、分離装置2-2は、XMLデータから要素の内容（テキスト）と構造（要素名および構造）を分離する。構文（値）生成装置2-3は、要素の構造からASN.1抽象構文（値）を

生成し、転送構文生成装置2-4は、ASN.1転送構文を生成する。圧縮装置2-5は分離したテキストを圧縮し、併合装置2-6は圧縮テキストとASN.1転送構文を併合して符号化されたXMLデータを生成する。

【0024】このような符号化システムを用いた本実施の形態の符号化方法の手順の概要は以下のとおりである。

（ステップ2-1） 構文（型）生成装置2-1を用いてDTDをASN.1抽象構文（型）に変換する。

10 （ステップ2-2） 分離装置2-2を用いてステップ2-1のDTDに従うXMLデータからテキストを分離する。なお、ステップ2-2は、ステップ2-1と並行に行われても良く、ステップ2-1よりも先に行われても良い。

（ステップ2-3） 構文（値）生成装置2-3を用いてステップ2-2でテキストが分離された後のXMLデータ（要素の構文）を、ステップ2-1のASN.1抽象構文（型）に従うASN.1抽象構文（値）に変換する。

20 （ステップ2-4） 転送構文生成装置2-4を用いてステップ2-3のASN.1抽象構文（値）をASN.1転送構文に変換する。

（ステップ2-5） 圧縮装置2-5を用いてステップ2-2で分離されたテキストを圧縮する。なお、ステップ2-5はステップ2-3～ステップ2-4と並行に行われる。

30 （ステップ2-6） 併合装置2-6を用いてステップ2-4で生成されたASN.1転送構文とステップ2-5で生成された圧縮テキストを併合し、符号化XMLデータを生成する。

【0025】なお、ASN.1抽象構文（値）をASN.1転送構文に変換するには、ASN.1で規定されている符号化規則に従う。そのような規則には、BERやDER、PERなどがある。特に、PERは、ASN.1抽象構文（型）から一意に決定される型やその値などを符号化しないため、符号化効率がよい。

【0026】以下、前記各ステップを詳細に説明する。

2.2 ステップ2-1

ステップ2-1では、DTDをASN.1抽象構文

40 （型）に変換する。以下では、内容モデルのパターンごとに、その方法を説明する。

【0027】2.2.1 要素内容(element content)
XMLの要素内容は、要素名とオペレータの組み合わせで構成される。通常要素名は文字列で表される。オペレータは要素内容における子要素の出現順と出現回数を指定するための演算子である。周知のとおりXMLでは、「」、「|」、「?」、「*」、「+」の各オペレータが許可される。原則として、要素名はASN.1抽象構文では識別子として、オペレータは同じくASN.1抽象構文では型として表現する。ただし、識別子はBE

Rなどでは符号化されないため、要素名を識別子とする必然性はない。以下では、オペレータごとに、それをどのような型で表現するかを説明する。

【0028】2. 2. 1. 1 「,」オペレータ
「,」オペレータは、ASN. 1抽象構文ではsequence型で表現する。例えば、DTDとして、
<!ELEMENT a (b,c)>
<!ELEMENT b (#PCDATA)>
<!ELEMENT c (#PCDATA)>

が与えられるとき、ASN. 1抽象構文(型)は、

A ::= SEQUENCE {

b B,
c C}
B ::= NULL
C ::= NULL

となる。ここで、A、B、Cは、便宜上導入した型参照である。型参照は、衝突さえしなければどのようなものでもよい。また、テキストはステップ2-2で分離されるため、BやCはnull型とする。

【0029】2. 2. 1. 2 「|」オペレータ
「|」オペレータは、ASN. 1抽象構文ではchoice型で表現する。例えば、DTDとして、
<!ELEMENT a (b|c)>
<!ELEMENT b (#PCDATA)>
<!ELEMENT c (#PCDATA)>

が与えられているとき、ASN. 1抽象構文(型)は、

A ::= CHOICE {
b B,
c C}
B ::= NULL
C ::= NULL

となる。

【0030】2. 2. 1. 3 「?」オペレータ
「?」オペレータは、ASN. 1抽象構文ではsequence型とキーワード「OPTIONAL」の組み合わせで表現する。例えば、DTDとして、
<!ELEMENT a (b?)>
<!ELEMENT b (#PCDATA)>

が与えらるとき、ASN. 1抽象構文(型)は、

A ::= SEQUENCE {
b B OPTIONAL}
B ::= NULL

となる。

【0031】2. 2. 1. 4 「*」オペレータ
「*」オペレータは、ASN. 1抽象構文ではsequence-of型で表現する。例えば、DTDとして、
<!ELEMENT a (b*)>
<!ELEMENT b (#PCDATA)>

が与えられるとき、ASN. 1抽象構文(型)は、

A ::= SEQUENCE OF B
B ::= NULL

となる。

【0032】2. 2. 1. 5 「+」オペレータ
「+」オペレータは、ASN. 1抽象構文ではサイズが制限されたsequence-of型で表現する。ただし、サイズ制限は符号化には影響しないため、「+」オペレータに対する符号は「*」オペレータに対するそれと同じものになる。例えば、DTDとして、

<!ELEMENT a (b+)>
<!ELEMENT b (#PCDATA)>

が与えられるとき、ASN. 1抽象構文(型)は、

A ::= SEQUENCE SIZE (1..MAX) OF B
B ::= NULL
となる。

【0033】2. 2. 1. 6 オペレータなし
要素内容には、オペレータの適用されていない要素名が1つだけ指定されることもある。そのような要素名は、ASN. 1抽象構文ではdefined型で表現する。例えば、DTDとして、

<!ELEMENT a (b)>
<!ELEMENT b (#PCDATA)>

が与えられるとき、ASN. 1抽象構文(型)は、

A ::= B
B ::= NULLとなる。

【0034】2. 2. 2 混在内容(mixed content)
混在内容は、キーワード「#PCDATA」と1つ以上の要素名を「|」オペレータで結合した後、「*」オペレータを適用したものである。そこで、ASN. 1抽象構文ではchoice型とsequence-of型の組み合わせで表現する。例えば、DTDとして、

<!ELEMENT a (#PCDATA|b)*>
<!ELEMENT b (#PCDATA)>

が与えらるとき、ASN. 1抽象構文(型)は、

A ::= SEQUENCE OF CHOICE {
txt NULL,
b B}
B ::= NULL

40 となる。ここで、txtは便宜上導入した識別子であり、混在内容に含まれるテキストに対応する。

【0035】2. 2. 3 空要素(EMPTY)
空要素は、ASN. 1抽象構文ではnull型で表現する。例えば、DTDが、

<!ELEMENT a EMPTY>

の場合はASN. 1抽象構文(型)は、

A ::= NULL
となる。

【0036】2. 2. 4 任意要素(ANY)

50 任意要素は、キーワード「#PCDATA」とDTDで

宣言されたすべての要素名から構成される混在内容と等価である。従って、任意要素のASN. 1抽象構文における表現は、混在内容のそれに帰着される。

【0037】2. 3 ステップ2-2

ステップ2-2では、ステップ2-1のDTDに従うXMLデータからテキストを分離する。例えば、XMLデータが、

```
(a)
  (b)10(b)
  (c)20(c)
(a)
```

の場合、要素bから「10」が、要素cから「20」が分離される。その結果、XMLデータ（要素名および構造）は次のようになる。

```
(a)
  (b /)
  (c /)
(a)
```

【0038】分離されたテキストを圧縮する方法は、例えば、XMLi11などのように、要素ごとにまとめた後、圧縮することができる。なお、その他の圧縮方法を適用してももちろん良い。

【0039】2. 4 ステップ2-3

ステップ2-3では、ステップ2-2のXMLデータを、ステップ2-1のASN. 1抽象構文（型）に従うASN. 1抽象構文（値）に変換する。例えば、ASN. 1抽象構文（型）が、

```
A ::= SEQUENCE {
  b B,
  c C }
B ::= NULL
C ::= NULL
```

の場合、以下のXMLデータ（要素名および構造）、

```
(a)
  (b /)
  (c /)
(a)
```

をASN. 1抽象構文（値）に変換すると、

```
a A ::= {
  b NULL,
  c NULL }
```

となる。

【0040】2. 5 ステップ2-4

ステップ2-4では、ASN. 1で規定されている符号化規則に従って、ステップ2-3のASN. 1抽象構文（値）をASN. 1転送構文に変換する。そのような規則にはBERやDER、PER（ALIGNED/UNALIGNED）などがある。しかしながら、符号化効率を向上する観点からPER（UNALIGNED）を

利用することが好ましい。尤も、BER、DER、PER（ALIGNED）を利用してもよいことは勿論である。

【0041】2. 6 ステップ2-5

ステップ2-5では、テキストを圧縮する。圧縮の具体的方法については、周知のLZ77等を用いる。その他、従来技術の項で説明した公知技術を用いることができる。圧縮は、要素ごとに行われても良く、また、各要素をまとめて圧縮しても良い。

10 【0042】2. 7 ステップ2-6

ステップ2-6では、ASN. 1転送構文と圧縮テキストを併合する。両データを単に結合することも可能であるが、復号時の分離を考慮したセパレータをデータ間に挿入し、あるいはデータビット数情報を持つヘッダ等を付加しても良い。

【0043】3. 復号化手順

3. 1 システム構成と手順の概要

図2は本実施の形態の復号化システムの一例をその機能について示したブロック図である。本実施の形態の復号化システムは、構文（型）生成装置3-1、転送構文復号装置3-2、抽象構文復号装置3-3、併合装置3-4、分離装置3-5、解凍装置3-6を有する。

【0044】構文（型）生成装置3-1は、前記した構文（型）生成装置2-1と同様に、DTDをASN. 1抽象構文（型）に変換する。分離装置3-5は、符号化されたXMLデータからASN. 1転送構文と圧縮テキストを分離し、転送構文復号装置3-2は、ASN. 1転送構文をASN. 1抽象構文（型）に従うASN. 1抽象構文（値）に変換する。抽象構文復号装置3-3

30 は、ASN. 1抽象構文（値）をDTDに従うXMLデータ（要素名と構造）に変換する。併合装置3-4は復号化されたテキスト（要素の内容）とXMLデータ（要素名と構造）を併合し、XMLデータを生成する。解凍装置3-6は、圧縮テキストを解凍する。

【0045】このような復号化システムを用いた本実施の形態の復号化方法の手順は、前記符号化の手順をほぼ逆に行う。その概要は以下のとおりである。

（ステップ3-1） DTDをASN. 1抽象構文（型）に変換する。

40 （ステップ3-2） 符号化XMLデータを圧縮テキストとASN. 1転送構文に分離する。なお、ステップ3-2は、ステップ3-1と並行に行われても良く、ステップ3-1よりも先に行われても良い。

（ステップ3-3） ASN. 1転送構文を、ステップ3-1のASN. 1抽象構文（型）に従うASN. 1抽象構文（値）に変換する。

（ステップ3-4） ステップ3-3のASN. 1抽象構文（値）を、ステップ3-1のDTDに従うXMLデータ（要素名と構造）に変換する。

50 （ステップ3-5） ステップ3-2で分離した圧縮デ

キストを解凍する。

(ステップ 3-6) ステップ 3-4 の XML データ (要素名と構造) にステップ 3-5 で解凍したテキストを結合する。

【0046】なお、前記ステップ 3-1 ～ステップ 3-6 における復号化の各処理は、それに相当する符号化の場合の逆であって自明である。よって、その詳細な説明を省略する。以下、DTD の変換およびその DTD に従う XML データの PER による符号化の具体例を示す。

【0047】4. 「,」オペレータを含む場合
ここでは、「,」オペレータを含む DTD と、それに従う XML データを例にとる。以下の DTD、

```
<!ELEMENT a (b,c)>
<!ELEMENT b (#PCDATA)>
<!ELEMENT c (#PCDATA)>
```

が与えられている場合、前記ステップ 2-1 で生成される ASN. 1 抽象構文 (型) は、

```
A ::= SEQUENCE {
    b B,
    c C }
B ::= NULL
C ::= NULL
```

となる。このような DTD に従う符号化対象の XML データとして、

```
<a>
    <b>10</b>
    <c>20</c>
</a>
```

を例示すれば、前記ステップ 2-2 で分離される要素 (XML データの要素名と構造) は、

```
<a>
    <b />
    <c />
</a>
```

となる。この要素から前記ステップ 2-3 で生成される ASN. 1 抽象構文 (値) は、

```
a A ::= {
    b NULL,
    c NULL }
```

となる。

【0048】PER では、sequence 型の値は、原則としてその構成要素の値がその順番で符号化される。しかし、ここでは b や c の値は null であり、null は空ビット列に符号化されるため、a の値は空ビット列となる。その場合に、符号は例外的に、00000000(2)

のようになる。このようにして、ASN. 1 転送構文が生成される。なお、このような ASN. 1 転送構文の生成は前記ステップ 2-4 で行われる。また、下付き文字

の“(2)”は2進数であることを示す。

【0049】その後、前記ステップ 2-5 で圧縮処理された圧縮テキストと ASN. 1 転送構文 (ここでは「00000000(2)」) が前記ステップ 2-6 で併合され、符号化 XML データが生成される。

【0050】5. 「|」オペレータを含む場合
ここでは、「|」オペレータを含む DTD と、それに従う XML データを例にとる。以下の DTD、

```
<!ELEMENT a (b|c)>
10 <!ELEMENT b (#PCDATA)>
<!ELEMENT c (#PCDATA)>
```

が与えられている場合、前記ステップ 2-1 で生成される ASN. 1 抽象構文 (型) は、

```
A ::= CHOICE {
    b B,
    c C }
B ::= NULL
C ::= NULL
```

となる。このような DTD に従う符号化対象の XML データとして、

```
<a>
    <b>10</b>
</a>
```

を例示すれば、前記ステップ 2-2 で分離される要素 (XML データの要素名と構造) は、

```
<a>
    <b />
</a>
```

となる。この要素から前記ステップ 2-3 で生成される ASN. 1 抽象構文 (値) は、

```
a A ::= b:NULL
```

となる。

【0051】PER では、choice 型の値は、まず選択された構成要素のインデックス (0 ベース) が符号化され、次にその構成要素の値が符号化される。ここでは b が選択されているため、インデックスは 0 である。また、b の値は null である。従って、a の値は、0XXXXXXX(2)

のようになる。ここで、X は、8 ビットの倍数にするために付加されたパディング・ビットを表わす。このようにして、ASN. 1 転送構文が生成される。なお、圧縮テキストの生成および圧縮テキストと ASN. 1 転送構文の併合は前記“4. 「,」オペレータを含む場合”と同様である。

【0052】6. 「?」オペレータを含む場合
ここでは、「?」オペレータを含む DTD と、それに従う XML データを例にとる。以下の DTD、

```
<!ELEMENT a (b?,c)>
<!ELEMENT b (#PCDATA)>
50 <!ELEMENT c (#PCDATA)>
```

19

が与えられている場合、前記ステップ 2-1 で生成される ASN. 1 抽象構文 (型) は、

```
A ::= SEQUENCE {
    id0 SEQUENCE {
        b B OPTIONAL },
        c C }
B ::= NULL
C ::= NULL
```

となる。このような DTD に従う符号化対象の XML データとして、

```
<a>
  <b>10</b>
  <c>20</c>
</a>
```

を例示すれば、前記ステップ 2-2 で分離される要素 (XML データの要素名と構造) は、

```
<a>
  <b />
  <c />
</a>
```

となる。この要素から前記ステップ 2-3 で生成される ASN. 1 抽象構文 (値) は、

```
a A ::= {
    id0 {
        b NULL }
        c NULL }
```

となる。

【0053】PERでは、sequence型の構成要素が少なくとも1つ「OPTIONAL」と指定されている場合に、各構成要素の値が符号化される前に、どの構成要素が存在しているかを表わすビット列が付加される。ビットは、構成要素が存在する場合に1、存在しない場合に0となる。ここではbだけが「OPTIONAL」と指定されており、それが存在しているため、まずビット列1が付加される。次にbとcの値が符号化されるが、それらは共にnullであるため、aの値は、1XXXXXX(2)

のようになる。このようにして、ASN. 1 転送構文が生成される。なお、圧縮テキストの生成および圧縮テキストとASN. 1 転送構文の併合は前記“4. 「,」オペレータを含む場”と同様である。

【0054】7. 「*」オペレータを含む場合
ここでは、「*」オペレータを含むDTDと、それに従うXMLデータを例にとる。以下のDTD、

```
<!ELEMENT a (b*)>
<!ELEMENT b (#PCDATA)>
```

が与えられている場合、前記ステップ 2-1 で生成される ASN. 1 抽象構文 (型) は、

```
A ::= SEQUENCE OF B
```

20

B ::= NULL

となる。このような DTD に従う符号化対象の XML データとして、

```
<a>
  <b>10</b>
  <b>20</b>
</a>
```

を例示すれば、前記ステップ 2-2 で分離される要素 (XML データの要素名と構造) は、

```
<a>
  <b />
  <b />
</a>
```

となる。この要素から前記ステップ 2-3 で生成される ASN. 1 抽象構文 (値) は、

```
a A ::= {
    NULL,
    NULL }
```

となる。

20 【0055】PERでは、sequence-of型の値は、まずその構成要素の数が符号化され、次に各構成要素の値がその順番で符号化される。ここでは構成要素の数は2であり、まずその値が符号化される。次に各構成要素の値が符号化されるが、それらはすべてnullであるため、aの値は、00000010(2)のようになる。このようにして、ASN. 1 転送構文が生成される。なお、圧縮テキストの生成および圧縮テキストとASN. 1 転送構文の併合は前記“4. 「,」オペレータを含む場合”と同様である。

30 【0056】8. 混在内容を含む場合

ここでは、混在内容を含むDTDと、それに従うXMLデータを例にとる。以下のDTD、

```
<!ELEMENT a (#PCDATA|b)*>
<!ELEMENT b (#PCDATA)>
```

が与えられている場合、前記ステップ 2-1 で生成される ASN. 1 抽象構文 (型) は、

```
A ::= SEQUENCE OF CHOICE {
    txt NULL,
    b NULL }
```

40 となる。このような DTD に従う符号化対象の XML データとして、

```
<a>
  xxx
  <b>10</b>
  <b>20</b>
</a>
```

を例示すれば、前記ステップ 2-2 で分離される要素 (XML データの要素名と構造) は、

21
 (a)
 <txt />

となる。

【0057】なお、上記において「xxx」はテキストの内容を示し、「txt」は要素がテキストであることを示す。この要素から前記ステップ2-3で生成されるASN. 1抽象構文(値)は、

a A ::= {
 txt:NULL,
 b:NULL,
 b:NULL }

となる。

【0058】sequence-of型の値とchoice型の値をPERで符号化する方法は、前述の通りである。ここでは、構成要素の数は3である。また、構成要素としてtxt、b、bがその順番で選択されているため、インデクスは0、1、1となる。それらの値は

00000011(2)

011XXXX(2)

のようになる。このようにして、ASN. 1転送構文が生成される。なお、圧縮テキストの生成および圧縮テキストとASN. 1転送構文の併合は前記“4. 「,」オペレータを含む場合」と同様である。

【0059】(実施の形態2)前記実施の形態1では、XMLデータは要素とテキストだけから構成され、その他のもの(例えば、属性や処理命令)は別の方法で管理されると仮定して説明を行った。本実施の形態では、属性や処理命令等が含まれる場合の処理の一例を説明する。すなわち、それら属性等を特別な要素で表現し、XMLデータの中に埋め込む方法を例示する。この場合、文法も変更する必要がある。なお、本実施の形態の例に関わらず、属性等をXMLデータから分離して、XPointerなどと共に別に格納する方法を採用することも可能である。このような属性等を特別な要素としてXMLデータに埋め込む方策は、前記実施の形態1の処理の前処理および後処理として把握することが可能である。これにより広い範囲のXMLデータを扱えるようになる。以下では、そのような前処理の例として、DTDで定義される項目をXMLデータの中に埋め込む方法を説明する。

【0060】9. 1 前処理システムと前処理の概要
 図3は、本実施の形態の前処理システムの一例をその機能について示したブロック図である。本実施の形態の前処理システムは、DTD変換装置9-1、XMLデータ変換装置9-2、符号化装置9-3を有する。

【0061】DTD変換装置9-1は、一般のDTDを

22

要素以外の項目についても考慮したDTD'に変換する。XMLデータ変換装置9-2は、一般のDTDに従うXMLデータをDTD'に従うXMLデータ'に変換する。符号化装置9-3は、実施の形態1の符号化システムである。

【0062】本実施の形態の前処理方法は、装置9-1によるDTDの変換および装置9-2によるXMLデータの変換の各ステップを有する。各々変換後、実施の形態1の符号化処理を行うことにより符号化XMLデータが生成される。

【0063】9. 2 DTDの変換処理

DTDには、要素型宣言以外に属性リスト宣言、エンティティ宣言、記法宣言の3つの宣言が含まれる。

【0064】エンティティ宣言で定義されるエンティティには、パース対象エンティティ、パース対象外エンティティ、パラメータ・エンティティがある。パース対象エンティティは、テキストや属性値など、あらゆる場所で参照される。このエンティティは、単純に展開する。また、パース対象外エンティティは、属性値でしか参照されえない。従って、属性を処理できれば十分である。一方、パラメータ・エンティティは、DTDの中でしか参照されえない。従って、パラメータ・エンティティは本実施の形態では考慮しない。

【0065】記法宣言で定義される記法は、属性値でしか参照されえない。従って、属性を処理できれば十分である。

【0066】以上の議論から、DTDに対する前処理は、要素型宣言と属性リスト宣言に対するそれに帰着される。

【0067】属性リスト宣言では、ある要素に付与される属性と、その属性の取りうる値が定義される。また、必要であれば、デフォルト値も定義される。属性は、「REQUIRED」であるか、「IMPLIED」であるか、デフォルト値が定義されているか、それは「FIXED」であるかで、要素に対する付与のされ方が変わる。従って、属性は原則として要素で表現することにするが、その定義に従って表現を多少変えることが好ましい。

【0068】属性リスト宣言に従って要素型宣言を変更する規則は、以下のとおりである。

(1) 属性は、要素(以下、属性要素)で表現する。要素名は、属性名から一意に決定できるようにする。例えば、“親の要素名”+“_”+“属性名”のようにする。ただし、要素名がすでに使用されていないことに注意する。以下では、この命名規則に従うことにする。

(2) 属性値はすべて「CDATA」として扱い、属性要素に含める。

(3) 属性要素は、親要素の子要素の先頭に挿入する。

(4) 「REQUIRED」属性は、要素で表現する。

(5) 「IMPLIED」属性は、「？」オペレータが

適用された要素で表現する。

(6) デフォルト値が定義されている属性は、「?」オペレータが適用された要素で表現する。

(7) 「FIXED」デフォルト値が定義されている属性は、無視する。

【0069】以下具体例をあげる。例えば、DTDとして、

```
<!ELEMENT a (b,c)>
<!ATTLIST a %CDATA %REQUIRED
          x ID %IMPLIED
          y (0|1|2) "0"
          z CDATA %FIXED "abc">
```

が与えられたとき、DTD' は、

```
<!ELEMENT a (a_w,a_x?,a_y?,b,c)>
<!ELEMENT a_w (#PCDATA)>
<!ELEMENT a_x (#PCDATA)>
<!ELEMENT a_y (#PCDATA)>
```

となる。

【0070】9.3 XMLデータの変換処理
ここでは、上記のDTDに従うXMLデータを上記DTD' に従うXMLデータ' に変換する。基本的には、要素に付与されている属性を属性要素に変換し、その属性要素をその要素の子要素の先頭に挿入すればよい。ただし、属性にデフォルト値が定義されており、属性値がそれと一致するのであれば、属性要素には変換しない。

【0071】例えば、9.2におけるDTDに従うXMLデータとして、

```
<a w="xyz" y="0" z="abc">
  <b>...</b>
  <c>...</c>
</a>
```

を例示すれば、XMLデータ' は、

```
<a>
  <a_wxyz</a_w>
  <b>...</b>
  <c>...</c>
</a>
```

に変換される。ここで、a要素にはx属性が付与されていないため、a_x要素は現れていない。また、y属性の値はデフォルト値と一致するため、a_y要素も現れていない。さらに、z属性は「FIXED」デフォルト値が定義されているため、削除されている。

【0072】9.4 後処理システムと後処理の概要
図4は、本実施の形態の後処理システムの一例をその機能について示したブロック図である。本実施の形態の後処理システムは、DTD変換装置9-4、復号化装置9-5、XMLデータ変換装置9-6を有する。

【0073】DTD変換装置9-4はDTD変換装置9-1と同様である。復号化装置9-5は、実施の形態1

の復号化システムである。XMLデータ変換装置9-6は、DTD' に従うXMLデータ' からDTDに従うXMLデータに変換する。

【0074】本実施の形態の後処理方法は、装置9-4によるDTDの変換および装置9-6によるXMLデータの変換の各ステップを有する。実施の形態1の復号化処理の後、これらステップを実行する。装置9-4によるDTDの変換は9.2と同様であり、装置9-6によるXMLデータの変換は9.3の変換を逆に行うことにより実行できる。よって、詳細な説明は省略する。なお、図中破線矢印は、必要があればDTDを参照することを意味する。属性要素を識別でき、かつその要素名から属性名を一意に決定できるのであれば、DTDを参照する必要はない。また、後処理は前処理がいかに行われたかに依存するので、前記前処理が異なる場合に本実施の形態の後処理もそれに併せて変更されることは勿論である。

【0075】10. 実施の形態の効果

上記した実施の形態1、2の符号化方法を用いれば、XMLデータを効率良く圧縮することが可能である。以下、XCompとの比較において、本実施の形態の効果を説明する。

【0076】XMLデータを符号化するとき、DTDから一意に決定される情報は符号化しないという点で、本手法とXCompは同じである。しかし、XCompは、XMLデータがある特定の構造をしている場合に、圧縮効率が悪くなる。具体的には、要素に「?」オペレータや「*」オペレータ（「+」オペレータを含む）が適用された場合である。

30 【0077】10.1 「?」オペレータ

本実施の形態の手法は要素が存在するかどうかをビット列で表現するのに対して、XCompは存在する要素（選択肢）に付与されたインデックスを羅列する。この差は、「?」オペレータが適用された要素がいくつか連続し、それらがすべて存在する場合に顕著に現れる。

【0078】例えば、DTD：

```
<!ELEMENT a (b?,c?)>
<!ELEMENT b (#PCDATA)>
<!ELEMENT c (#PCDATA)>
```

40 に従うXMLデータ：

```
<a>
  <b>xxx</b>
  <c>yyy</c>
</a>
```

を考える。このXMLデータには、要素bと要素cが共に存在している。このXMLデータは、本手法では11(2)に符号化される。ただし、簡単のため、パディング・ビットは付加していない。

【0079】一方、XCompでは、オートマトン：

	25			
	S1	S2	S3	S4
S1	-	b/1	c/2	$\epsilon/3$
S2	-	-	c/1	$\epsilon/2$
S3	-	-	-	ϵ/ϵ
S4	-	-	-	-

が作成されたとすると、インデクス列 11 が得られる。
このインデクス列を 0 ベースに変換し、必要最小限のビット数で符号化すると、符号 000 (2) が得られる。
本手法による符号と XComp によるそれを比較すると、本手法のほうが 1 ビットだけ短く符号化できることが分かる。

【0080】表 1 に、「？」オペレータが適用された要

	XComp(従来技術)			本発明
「？」の数	選択肢数	ビット数	トータルビット数	ビット数
1	2	1	1	1
2	3、2	2、1	3	2
3	4、3、2	2、2、1	5	3
4	5、4、3、2	3、2、2、1	8	4
5	6、5、4、3、2	3、3、2、2、1	11	5

【0082】10. 2 「*」オペレータ

本実施の形態の手法は存在する要素の数を符号化するのにに対して、XComp は存在する要素（選択肢）に付与されたインデクスを羅列する。この差は、「*」オペレータが適用された要素がいくつも存在する場合に顕著に現れる。

【0083】例えば、DTD：

<!ELEMENT a (b*)>

<!ELEMENT b (#PCDATA)>

に従う XML データ：

(a)

(b)xxx(/b)

(b)yyy(/b)

...

(b)zzz(/b)

(/a)

を考える。ここで、要素 b は 8 回出現しているとする。

この XML データは、本手法では 00001000

(2) に符号化される。

【0084】一方、XComp では、オートマトン：

	S0	S1
S0	b/1	$\epsilon/2$
S1	-	-

が作成されたとすると、インデクス列 11111111 2 が得られる。このインデクス列を 0 ベースに変換し、必要最小限のビット数で符号化すると、符号 000000001 (2) が得られる。本手法による符号と XComp によるそれを比較すると、本手法のほうが 1 ビットだけ短く符号化できることが分かる。

【0085】図 5 に、「*」オペレータが適用された要

26

素がいくつか連続し、それらがすべて存在する場合に、それらを符号化するのに必要なビット数を、本手法と XComp で比較したものを示す。表 1 から、本手法は XComp と同じか、それよりも効率よく符号化できることが分かる。一般には、「？」オペレータが適用された要素が n 個連続する場合に、XComp は $O(n \log n)$ ビットを必要とするのに対して、本手法は $O(n)$ ビットで十分である。

【0081】

【表 1】

20 素がいくつも存在する場合に、それらを符号化するのに必要なビット数を、本手法と XComp で比較したものを示す。同図から、要素が 8 つ以上存在する場合に、本手法は XComp よりも効率よく符号化できることが分かる。本手法は、一般には XComp と同じオーダのビット数を必要とするが、実質的にはそれよりも少ないビット数で十分である。

【0086】以上、本発明者によってなされた発明を発明の実施の形態に基づき具体的に説明したが、本発明は前記実施の形態に限定されるものではなく、その要旨を

30 逸脱しない範囲で種々変更可能である。

【0087】たとえば、本発明の手法は、複数のオペレータを組み合わせて適用しても良い。内容モデルで使用されるオペレータには、「|」、「?」、「*」（「+」を含む）の 4 つがある。従って、適用順序を考慮に入れると、それらの組み合わせは 16 通りある。各組み合わせに対して、前記実施の形態と同様に XML データを ASN. 1 抽象構文に変換することができる。よって、本発明を利用してあらゆるオペレータの組合せを含む XML データを符号化できる。

40 【0088】また、前記実施の形態で説明した手法のポイントの 1 つは、文法やそれに従う XML データの ASN. 1 抽象構文における表現を定義したことである。しかし、それとは別の表現も考えられる。

【0089】たとえば、「？」オペレータは、前記実施の形態では sequence 型とキーワード「OPTIONAL」の組み合わせで表現した。具体的には、DTD：

<!ELEMENT a (b?,c)>

は、ASN. 1 抽象構文：

27
 A ::= SEQUENCE {
 ido SEQUENCE {
 b B OPTIONAL },
 c C }

に変換される。

【0090】しかし、先のDTDはASN. 1抽象構文：

A ::= SEQUENCE {
 b B OPTIONAL,
 c C }

に変換されてもよい。このように変換されるほうが、BERやDERで符号化する場合に符号が短くなるため都合がよい。

【0091】ただし、この方法では、以下のDTD：

<!ELEMENT a (b?)>

<!ELEMENT a (b?|c)>

のような場合にうまく機能しないので、実施の形態のように表現するほうが妥当である。

【0092】また、前記実施の形態では、オペレータは何からの型で表現される。例えば、「,」オペレータは 20

A ::= SEQUENCE {
 b B }

と表現できる。

【0094】この表現法では、sequence型で統一的に表現することで、見た目には分かりやすくなる。また、オペレータをsequence型でラップすることに相当するので、複数のオペレータを組み合わせてもうまく機能することが分かる。ただし、BERやDERでは常に型が符号化されるため、型を挿入すればするほど符号は長くなる。従って、必要以上の型を挿入すべきではなく、その理由から前記実施の形態の手法で表現するほうが妥当である。

【0095】

【発明の効果】本願で開示される発明のうち、代表的なものによって得られる効果は、以下の通りである。すなわち、XMLデータの符号化（圧縮）効率を高めることが可能になり、また、属性等要素以外の記述を含めたXMLデータの符号化（圧縮）が可能になる。これにより、XMLデータの転送による通信負荷を軽減し、また、XMLデータのストレージ容量を少なくできる。

【図面の簡単な説明】

【図1】本発明の一実施の形態である符号化システムの

28

sequence型で表現され、「|」オペレータはchoice型で表現される。一方、オペレータの適用をプロダクションと見なし、すべてsequence型で表現するという方法も採れる。その場合に、sequence型で表現されていない「|」オペレータ、「*」オペレータ（「+」オペレータを含む）およびオペレータなしは、次のように表現できる。

【0093】すなわち、<!ELEMENT a(b|c)>は、

A ::= SEQUENCE {
 ido CHOICE {
 b B,
 c C } }

と表現できる。また、<!ELEMENT a (b*)>は、

A ::= SEQUENCE {
 b SEQUENCE OF B }

と表現できる。また、<!ELEMENT a (b)>は、

一例をその機能について示したブロック図である。

【図2】本発明の一実施の形態である復号化システムの一例をその機能について示したブロック図である。

【図3】本発明の一実施の形態である前処理システムの一例をその機能について示したブロック図である。

【図4】本発明の一実施の形態である後処理システムの一例をその機能について示したブロック図である。

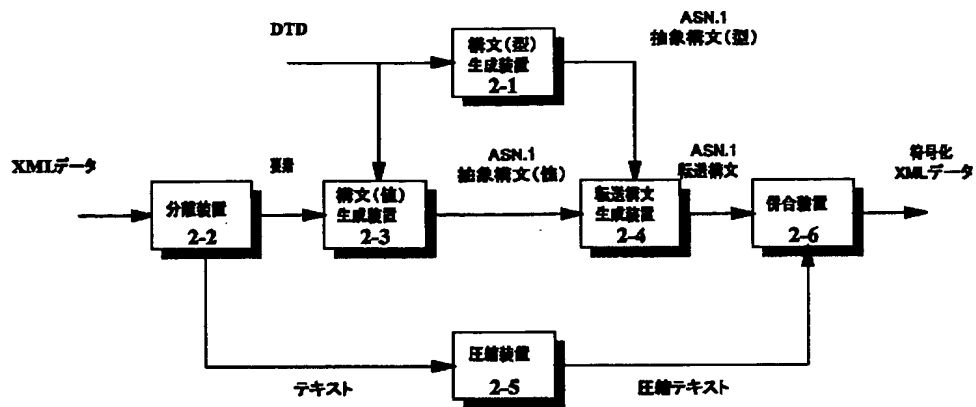
30 【図5】本実施の形態の符号化手法の符号化率をXCompとの比較において示した図である。

【符号の説明】

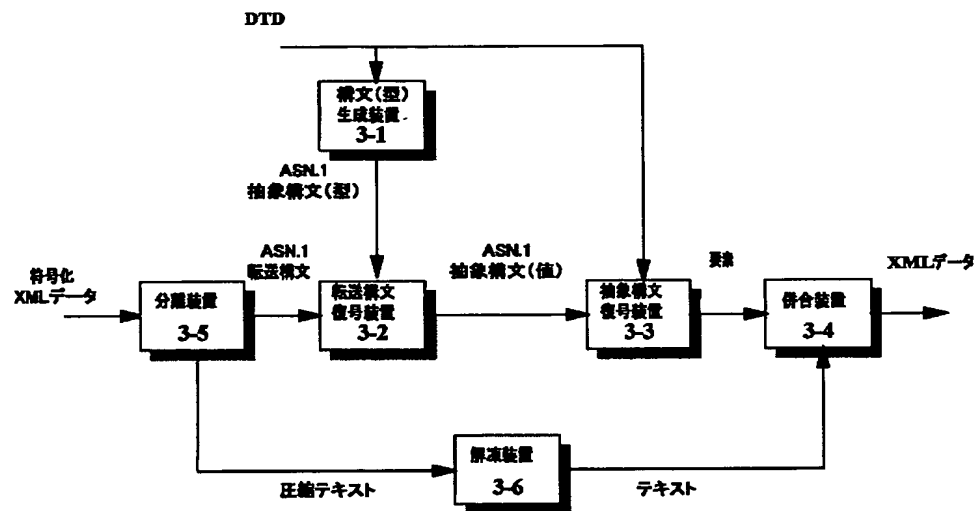
2-1…構文（型）生成装置、2-2…分離装置、2-3…構文（値）生成装置、2-4…転送構文生成装置、2-5…圧縮装置、2-6…併合装置、3-1…構文（型）生成装置、3-2…転送構文復号装置、3-3…抽象構文復号装置、3-4…併合装置、3-5…分離装置、3-6…解凍装置、9-1…DTD変換装置、9-2…XMLデータ変換装置、9-3…符号化装置、9-4…DTD変換装置、9-5…復号化装置、9-6…XMLデータ変換装置。

40

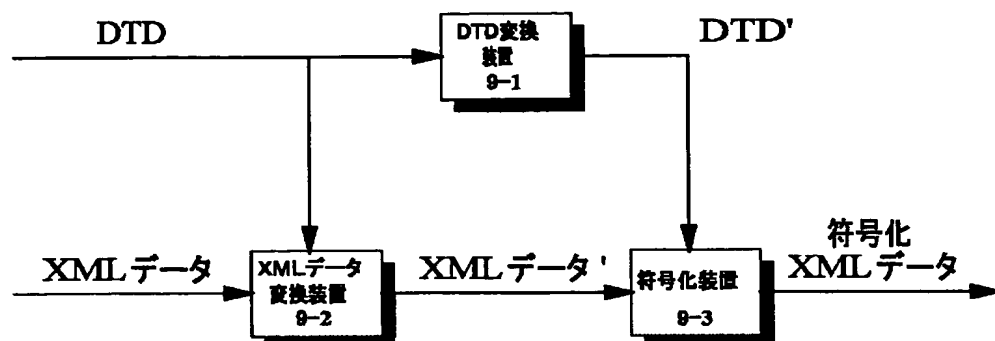
【図1】



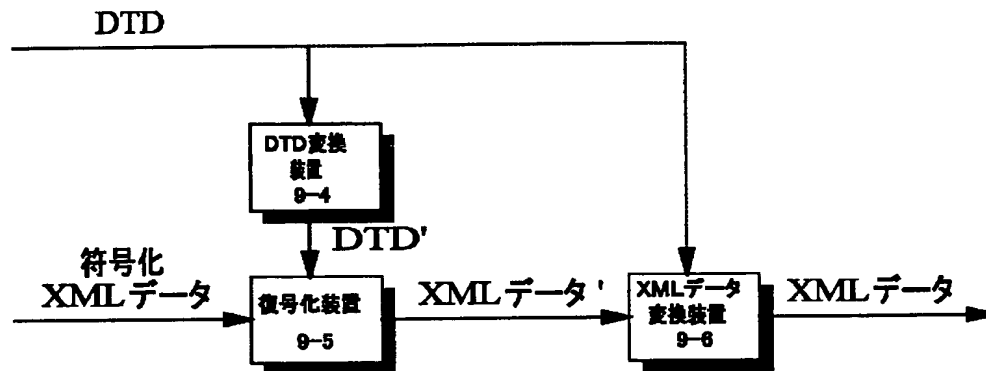
【図2】



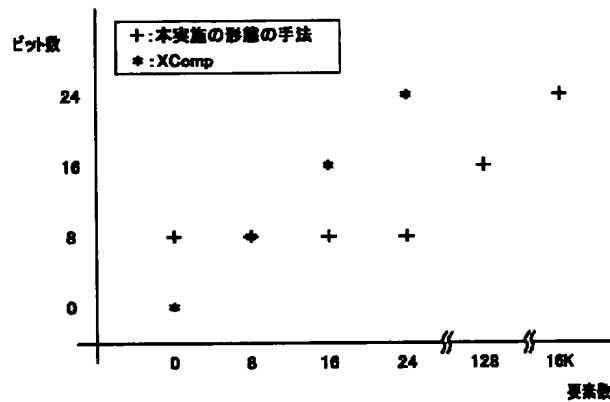
【図3】



【図 4】



【図 5】



フロントページの続き

(51) Int. Cl.⁷

識別記号

F I

テーマコード* (参考)

// G 0 6 F 17/21

5 7 0

G 0 6 F 17/21

5 7 0 G

(72) 発明者 今村 剛

神奈川県大和市下鶴間1623番地14 日本ア

イ・ピー・エム株式会社 東京基礎研究所

内

F ターム(参考) 5B009 SA07

5B082 GA01

5B089 HB10 KA08 KH04 KH28

5J064 AA02 BA11 BC02 BD03